

UNIVERSITY OF TECHNOLOGY SYDNEY
Faculty of Engineering and Information Technology

Transformer-Based Siamese Textual Encoder for Knowledge Base Completion

by

Mengyao Li

A THESIS SUBMITTED
IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE

Master of Analytics (Reasarch)

Sydney, Australia

2020

Certificate of Authorship/Originality

I, Mengyao Li, declare that this thesis, is submitted in fulfilment of the requirements for the award of Master of Analytics (Research), in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise reference or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Production Note:
Signature: Signature removed
prior to publication.

Date: 23/02/2020

Dedication

To my parents Xiuhua Ren and Guanhai Li

Acknowledgements

First of all, many thanks to my principal supervisor Dr. Jing Jiang and co-supervisor Dr. Guodong Long. During chasing the degree of research master, my supervisors gave me careful guidance and kind care in my study, research, and life. Their profound knowledge, rigorous academic attitude, and tireless study spirit made me admire, and left an indelible impression on me, which greatly influenced and inspired me. I will bear their edification in my mind, and these will definitely benefit me for my whole life.

Secondly, I want to thank my group mates, Bo Wang, Tao Shen, Zonghan Wu, Yang Li, Hao Huang, Peng Zhang for their help in all aspects of the model implementing and thesis writing. They gave me tremendous help with the questions I raised and put forward many valuable suggestions for my codes and thesis, which benefit me a lot.

Finally, I am deeply grateful to my parents, Xiuhua Ren and Guanhai Li. Every progress I have made is condensed by their selfless dedication, support, and encouragement; thanks to every teacher who has taught me in the growth of the patient teaching; thanks to all my friends who mutually helped to grow up together and the friendship deserves my cherish. These lovely people will continue to encourage me to move forward for more solid and comprehensive works in the future.

Mengyao Li

Sydney, Australia, 2020.

List of Publications

Journal Papers

- J-1. **Mengyao Li** and Jing Jiang, “Siamese Network for Knowledge Base Completion,” *Neural Processing Letters*, pre-acceptance, 2020.

Contents

Certificate	ii
Dedication	iii
Acknowledgments	iv
List of Publications	v
List of Figures	viii
Abstract	x
Abbreviation	xii
Notation	xiv
1 Introduction	1
1.1 Research Objectives	6
2 Literature Review	7
2.1 Graph Embedding Approaches for Link Prediction	7
2.1.1 Translation-based Embedding Approaches	8
2.1.2 Conv-based Embedding Approaches	12
2.1.3 GCN-based Embedding Approaches	14
2.2 Neural Textual Encoding in Natural Language Processing	16
2.2.1 Textual Encoding Approach	16
2.2.2 Siamese Neural Network	24
2.3 Textual Encoding Approach for Link Prediction	31

2.3.1	Individual Encoding Approach for KG	32
2.3.2	Contextualized Encoding Approach for KG	36
3	Methodology	41
3.1	Textual Encoding Approach for Knowledge Base Completion	41
3.2	Pre-Trained Transformer Encoder	42
3.3	Baseline Model: KG-BERT	44
3.4	Siamese Encoder for Knowledge Base Completion	45
3.5	Triple-Based Learning Objective	47
3.6	Model Efficiency during Inference	47
4	Experimental Results	49
4.1	Knowledge Graph	49
4.1.1	WordNet: A Large Lexical Graph of English	49
4.1.2	Freebase: Large-Scale Graph with Curated Knowledge	51
4.2	Experimental Setups	52
4.3	Evaluation Results on Link Prediction	54
4.4	Comparison to Baseline	54
4.5	Improvement Analysis	56
4.6	Error Analysis	58
5	Summary	60

List of Figures

2.1	One-dimensional Convolutional Neural Network (1D-CNN).	17
2.2	Basic Recurrent Neural Network (RNN) and its unroll status.	18
2.3	Scale dot-product and multi-head attention mechanisms.	21
2.4	The Transformer encoder.	22
2.5	Pre-training and fine-tuning procedures for BERT.	24
2.6	Siamese architecture for image classification.	26
2.7	Fully-convolutional Siamese architecture.	27
2.8	Sentence-BERT for both classification (left) and regression (right)).	29
2.9	An illustration of training representations for entities in a KB.	32
2.10	An illustration of neural tensor network (NTN).	33
2.11	Convolutional neural network encoder for plausibility score of a triple.	34
2.12	Comparison between TransE model (left) and SSP model (right) on a toy example.	35
2.13	Model structure of COMmonsEnse Transformers (COMET).	38
2.14	Masking Strategy for COMmonsEnse Transformers (COMET).	39
3.1	Bidirectional Encoder Representations from Transformers (BERT) in pre-training on two self-supervised tasks (left) and fine-tuning on various natural language processing tasks (right).	42
3.2	KG-BERT model that applies pre-trained Transformer encoder to knowledge base completion tasks.	44

3.3	The proposed Siamese encoder for knowledge base completion.	46
-----	---	----

ABSTRACT

Transformer-Based Siamese Textual Encoder for Knowledge Base Completion

by

Mengyao Li

Knowledge graph is constructed as a directed graph and regarded as a kind of knowledge base. Its nodes and edges represent named entities and relations between the entities, respectively. Knowledge graph is normally denoted as a collection of triples, i.e., (head/subject entity, relation, tail/object entity), for feasible processing. However, incompleteness or sparsity issue usually exists in knowledge graph since the knowledge graph are constructed via human labeling. Therefore, *link prediction* task is proposed to complete a knowledge graph via predicting the missing links.

Previous approaches for link prediction on knowledge graph are mainly based on graph embedding, such as translation-based embedding approaches, convolutional-based embedding approaches, and graph convolutional network (GCN-based) approaches. These traditional graph embedding-based approaches straightforwardly learn the embeddings by considering a knowledge base's structure but are inherently vulnerable to the graph's sparseness or incompleteness issue. In contrast, some recent approaches, named textual encoding approaches here, learn the embeddings from a natural language processing perspective. That is, except for the structural information in a knowledge graph, the semantic information of the entities and relations, like their texts, mentions, and descriptions, are taken into account to enrich the embeddings of the knowledge graph. Hence, they capture such structured knowledge from the semantic perspective and employ deep neural text encoder to model graph triples in semantic space, but fail to trade off the contextual features or model's efficiency.

Therefore, in line with previous textual encoding approaches that represent graph components like entities or triples into latent semantic space, we propose a Siamese textual encoder operating on each graph triple from knowledge base, where the contextual features between a head/tail entity and a relation are well-captured to highlight relation-aware entity embedding while a Siamese structure is also adapted to avoid combinatorial explosion during inference of link prediction.

In the experiments, the proposed approach reaches the best or comparable performance on two link prediction datasets. And further compared to its baseline (i.e., a state-of-the-art textual encoding approach), our approach can accelerate the inference procedure by one or two orders of magnitude with an even better quality of predictions.

Abbreviation

KG - Knowledge Graph

KB - Knowledge Base

Conv - Convolution

CNN - Convolutional Neural Network

GCN - Graph Convolutional Network

R-GCN - Relational Graph Convolutional Network

NLP - natural language processing

LSTM - Long Short-Term Memory

GRU - Gated Recurrent Unit

RNN - Recurrent Neural Network

BPE - byte-pair-encoding

CBoW - Continuous bag-of-words

FFN - feed-forward network

SAN - self-attention network

LM - language model

MLM - masked language model

NSP - next sentence prediction

ELMo - Embeddings from Language Models

BERT - Bidirectional Encoder Representations from Transformers

RoBERTa - Robustly Optimized BERT Pretraining Approach

KG-BERT - Knowledge Graph BERT

GPT - Generative Pre-Training

SSP - semantic space projection

PMI - point-wise mutual information

Nomenclature and Notation

Numbers and Arrays

a	A scalar (integer or real)
\boldsymbol{a}	A vector
\boldsymbol{A}	A matrix
\mathbf{A}	A tensor
\boldsymbol{I}_n	Identity matrix with n rows and n columns
\boldsymbol{I}	Identity matrix with dimensionality implied by context
$\boldsymbol{e}^{(i)}$	Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at position i
$\text{diag}(\boldsymbol{a})$	A square, diagonal matrix with diagonal entries given by \boldsymbol{a}
a	A scalar random variable
\boldsymbol{a}	A vector-valued random variable
\boldsymbol{A}	A matrix-valued random variable

Sets and Graphs

\mathbb{A}	A set
\mathbb{R}	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and n
$[a, b]$	The real interval including a and b
$(a, b]$	The real interval excluding a but including b
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of \mathbb{A} that are not in \mathbb{B}

\mathcal{G}	A graph
$Pa_{\mathcal{G}}(\mathbf{x}_i)$	The parents of \mathbf{x}_i in \mathcal{G}

Indexing

a_i	Element i of vector \mathbf{a} , with indexing starting at 1
a_{-i}	All elements of vector \mathbf{a} except for element i
$A_{i,j}$	Element i, j of matrix \mathbf{A}
$\mathbf{A}_{i,:}$	Row i of matrix \mathbf{A}
$\mathbf{A}_{:,i}$	Column i of matrix \mathbf{A}
$A_{i,j,k}$	Element (i, j, k) of a 3-D tensor \mathbf{A}
$\mathbf{A}_{::,i}$	2-D slice of a 3-D tensor
a_i	Element i of the random vector \mathbf{a}

Calculus

$\frac{dy}{dx}$	Derivative of y with respect to x
$\frac{\partial y}{\partial x}$	Partial derivative of y with respect to x
$\nabla_{\mathbf{x}} y$	Gradient of y with respect to \mathbf{x}
$\nabla_{\mathbf{X}} y$	Matrix derivatives of y with respect to \mathbf{X}
$\nabla_{\mathbf{X}} y$	Tensor containing derivatives of y with respect to \mathbf{X}
$\frac{\partial f}{\partial \mathbf{x}}$	Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or $\mathbf{H}(f)(\mathbf{x})$	The Hessian matrix of f at input point \mathbf{x}
$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of \mathbf{x}
$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to \mathbf{x} over the set \mathbb{S}

Probability and Information Theory

$P(\mathbf{a})$	A probability distribution over a discrete variable
$p(\mathbf{a})$	A probability distribution over a continuous variable, or over a variable whose type has not been specified
$\mathbf{a} \sim P$	Random variable \mathbf{a} has distribution P
$\mathbb{E}_{\mathbf{x} \sim P}[f(\mathbf{x})]$ or $\mathbb{E}f(\mathbf{x})$	Expectation of $f(\mathbf{x})$ with respect to $P(\mathbf{x})$
$\text{Var}(f(\mathbf{x}))$	Variance of $f(\mathbf{x})$ under $P(\mathbf{x})$
$\text{Cov}(f(\mathbf{x}), g(\mathbf{x}))$	Covariance of $f(\mathbf{x})$ and $g(\mathbf{x})$ under $P(\mathbf{x})$
$H(\mathbf{x})$	Shannon entropy of the random variable \mathbf{x}

$D_{\text{KL}}(P\ Q)$	Kullback-Leibler divergence of P and Q
$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian distribution over \mathbf{x} with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$

Functions

$f : \mathbb{A} \rightarrow \mathbb{B}$	The function f with domain \mathbb{A} and range \mathbb{B}
$f \circ g$	Composition of the functions f and g
$f(\mathbf{x}; \boldsymbol{\theta})$	A function of \mathbf{x} parametrized by $\boldsymbol{\theta}$. (Sometimes we write $f(\mathbf{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation)
$\log x$	Natural logarithm of x
$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
$\ \mathbf{x}\ _p$	L^p norm of \mathbf{x}
$\ \mathbf{x}\ $	L^2 norm of \mathbf{x}
x^+	Positive part of x , i.e., $\max(0, x)$
$\mathbf{1}_{\text{condition}}$	is 1 if the condition is true, 0 otherwise